

AD \_\_\_\_\_

Award Number: DAMD17-03-1-0697

TITLE: Computerized Identification of Normal Mammograms

PRINCIPAL INVESTIGATOR: Robert M. Nishikawa, Ph.D.

CONTRACTING ORGANIZATION: Chicago University  
Chicago, Illinois 60637

REPORT DATE: October 2004

TYPE OF REPORT: Annual

PREPARED FOR: U.S. Army Medical Research and Materiel Command  
Fort Detrick, Maryland 21702-5012

DISTRIBUTION STATEMENT: Approved for Public Release;  
Distribution Unlimited

The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision unless so designated by other documentation.

20050603 196

**REPORT DOCUMENTATION PAGE**Form Approved  
OMB No. 074-0188

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing this collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503

<b>1. AGENCY USE ONLY</b> (Leave blank)		<b>2. REPORT DATE</b> October 2004	<b>3. REPORT TYPE AND DATES COVERED</b> Annual (30 Sep 2003 - 29 Sep 2004)	
<b>4. TITLE AND SUBTITLE</b> Computerized Identification of Normal Mammograms			<b>5. FUNDING NUMBERS</b> DAMD17-03-1-0697	
<b>6. AUTHOR(S)</b> Robert M. Nishikawa, Ph.D.				
<b>7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)</b> Chicago University Chicago, Illinois 60637  <b>E-Mail:</b> r-nishikawa@uchicago.edu			<b>8. PERFORMING ORGANIZATION REPORT NUMBER</b>	
<b>9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES)</b> U.S. Army Medical Research and Materiel Command Fort Detrick, Maryland 21702-5012			<b>10. SPONSORING / MONITORING AGENCY REPORT NUMBER</b>	
<b>11. SUPPLEMENTARY NOTES</b>				
<b>12a. DISTRIBUTION / AVAILABILITY STATEMENT</b> Approved for Public Release; Distribution Unlimited				<b>12b. DISTRIBUTION CODE</b>
<b>13. ABSTRACT (Maximum 200 Words)</b> <p>The purpose of this concept award project is to develop an automated method to identify normal mammograms, that is those without breast disease. This is a new paradigm in computer-aided diagnosis (CAD), since all other CAD schemes identify breast cancer. We are relying on the natural pattern of glandular tissue in the normal breast, which radiates out from the nipple. Breast cancer disturbs this pattern. We have developed a database of 3000 regions of interest (ROIs) of normal breast tissue and 200 regions containing a portion of a breast cancer. Each region was automatically extracted from a mammogram that was reduced in size and preprocessed using a wavelet filter. We are using these ROIs to train an artificial neural network called a self-organizing map (SOM) to learn the mammographic pattern of normal breast tissue. SOM are self-learning classifiers that categorize input data into a user-defined number of distinct classes. To date, we have been unsuccessful in training the SOM to categorize normal and abnormal ROIs in a reliable manner. We are in the process of changing our training protocol.</p>				
<b>14. SUBJECT TERMS</b> Mammography, detection, computer-aided detection, normal				<b>15. NUMBER OF PAGES</b> 10
				<b>16. PRICE CODE</b>
<b>17. SECURITY CLASSIFICATION OF REPORT</b> Unclassified	<b>18. SECURITY CLASSIFICATION OF THIS PAGE</b> Unclassified	<b>19. SECURITY CLASSIFICATION OF ABSTRACT</b> Unclassified	<b>20. LIMITATION OF ABSTRACT</b> Unlimited	

## Table of Contents

Cover.....	1
SF 298.....	2
Table of Contents.....	3
Introduction.....	4
Body.....	4
Key Research Accomplishments.....	9
Reportable Outcomes.....	9
Conclusions.....	9
References.....	10
Appendices.....	

## 4. INTRODUCTION

Computer-aided detection (CAdE) systems have sensitivities at least equal to radiologists, 80-90% depending on the system, but the false detection rate is more than a magnitude higher than that of radiologists (on average the computer has 2 false detections per case, whereas a radiologist will have a one false positive every 10 cases). Because of the high false detection rate, a radiologist must review virtually every mammogram. Instead of locating abnormalities in mammograms, as is done all current CAdE systems, we propose to develop a method for determining normal mammograms. Initially, our approach would allow the radiologist to read only those cases that are judged to be not normal, reducing the number of cases reviewed potentially to 90% or better, allowing for more time to read cases that are more likely to contain a malignancy. Ultimately, if our approach is effective and optimized, it could be used as a front-end (triage system) to conventional CAD schemes that could be optimized to run on the "not normal" cases. Furthermore, we believe that the ultimate performance of CAD systems will not improve to the level of a radiologist using the current paradigm. A normal breast has a pattern of structures radiating out from the nipple. A cancer can disrupt this pattern. Our approach is to use this radiating pattern as a basis for recognizing normal mammograms. We will process the image to highlight the radiating pattern. Then by taking small regions of interest, we will train a classifier to recognize normal ROIs. The classifier used in this study is a specialized artificial neural network called a self-organizing map (SOM) {1}.

## 5. BODY

### 5.1. Tasks

*Task 1.* Process image to highlight ductal system

- a. Assemble 2,000 consecutive digitized normal screening exams and 100 cancer exams (cc views only) from an existing database of 25,000 consecutive screening mammograms.
- b. Create 3 datasets: (i) development set (500 normals); (ii) training set (1000 normals and 75 cancers); and (iii) testing set (500 normals and 25 cancers).
- c. Reduce image size by a factor of 10, testing different methods such as mean, maximum, median, and rank order.
- d. Implement two processing techniques, morphological operators and a linear detection algorithm developed by Zwiggelaar *et al.* (using development dataset)

*Task 2.* Train support vector machine to recognize normal mammogram:

- a. Train support vector machine (using training dataset)
- b. Measure the performance of the technique (using testing dataset)

#### 5.1.a Assemble databases

In a previous 5-year project, we digitized over 20,000 consecutive screen-film mammograms to 10 bits and 100-micron pixel size {2}. From this dataset, we have assembled 54 cancer-free consecutive cases and 5 cancer cases, collecting only the cranio-caudal (cc) views.

The abnormal cases contain a mass that was biopsied and found to be malignant. The normal cases were obtained by reading all the radiology reports for that patient. In a separate process, these reports had all patient identifiers removed and all reports from a single patient were placed in a single file and identified by the study number that was generated previously to allow the radiology report to be associated with the image. The study number is not traceable to any patient identifier. The mammograms are devoid of patient identifiers. To be considered normal, the case must have had at least a two-year period in which the mammograms were considered normal. Further, we selected from these cases, cases that were free of any type of lesion, including obvious benign findings such as lymph nodes and calcified vessels. This subset was used in the development data set.

Two other datasets are being created: a training set and a testing set. The exact composition of those data sets still needs to be decided. We are uncertain at this time whether to include, for example, obvious benign findings. If the SOM works well it may be able to classify obvious benign findings separate from suspicious lesions. This will need to be evaluated during the training phase of the study.

In the development phase, the goal is to understand how to pre-process the image and to understand how the SOM works. To do this, we need only a small database with very few cancer cases, in part because a large number of regions-of-interest (ROIs) can be selected from each image. As part of the development phase, we will be able to determine the number of cases needed to adequately train and test the SOM. Therefore, we have not finished collecting cases. We do have 300 normal and 70 cancer cases already selected for use in the training and testing data sets. If necessary, we will collect more cases from the 20,000 already digitized cases.

#### 5.1.b. Preprocess the mammograms

The 54 normal cases and the 5 abnormal cases were preprocessed to produce ROIs the either contain a portion of a cancer or are cancer free. This was done in four steps.

Step 1. The breast border was determined using software previously developed in our laboratory {3}.

Step 2. Wavelet decomposition was applied to the image using a bi-orthogonal spline mother wavelet implemented in MATLAB. All mother wavelets available in MATLAB were tested, but the bi-orthogonal spline gave the best visual result. This mother wavelet was used by Strickland in his study of detecting mammographic calcifications using wavelets {4}. We constructed the magnitude image from the horizontal and vertical components of the wavelet transform using level 3 (see Fig. 1). We originally had planned to implement a morphological operator and a linear detection algorithm developed by Zwiggelaar *et al.* {5} We spent several weeks implementing the Zwiggelaar method but it did not produce satisfactory results. To save time, we implemented the wavelet filtering method in MATLAB.

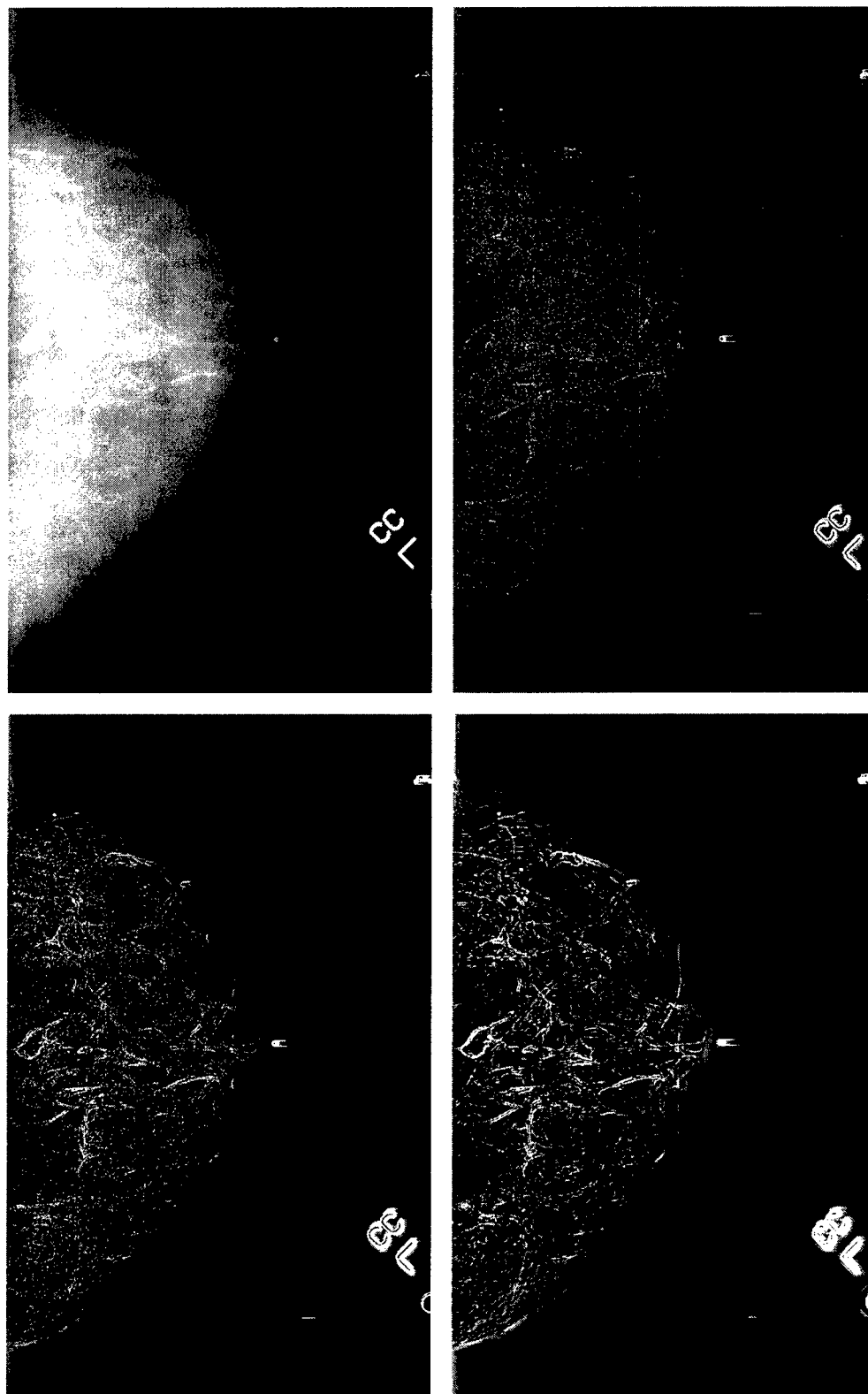


Figure 1. Illustration of the wavelet preprocessing. The original image is shown in the upper left. The other three images are the magnitude image of the wavelet transform for level 1 (upper right), level 2 (lower left) and level 3 (lower right). We used level 3 in this study.

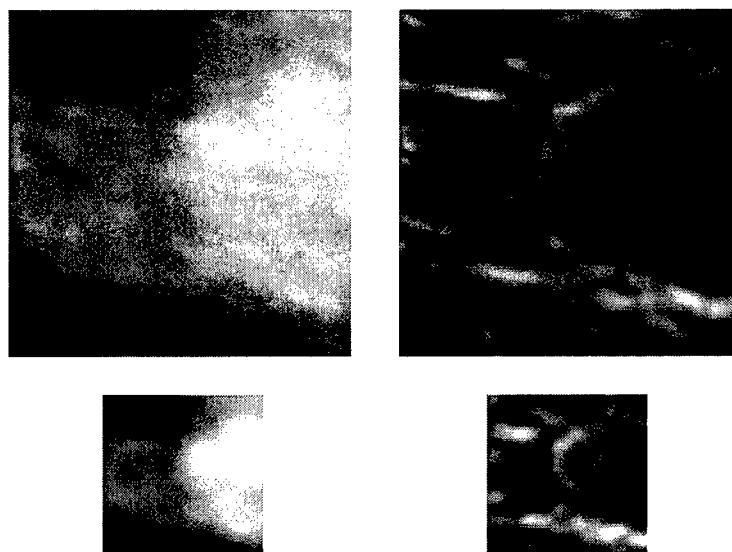


Figure 2. An illustration of the down sampling of the regions-of-interest (ROI). The top row show a 128x128 ROI extracted from the original image (left) and the wavelet processed image (right). The bottom row shows the two images after 8x8 pixel averaging. These two ROIs have been enlarged by a factor of 4. The ROI on the bottom right is representative of the ROIs used to train the SOM. The image on the bottom left is shown only for comparison purposes.

3. Based on the estimated breast border the largest rectangle that fit in the breast boundary was extracted from the wavelet image. From this rectangle, overlapping candidate ROIs that were 128x128 pixels in size were extracted. Each candidate ROI was shifted by 64 pixels from the previous candidate ROI. For each candidate ROI, a histogram of its pixel values was calculated. An upper and lower bound threshold was used to filter out "partial ROIs" (i.e., those that include non-breast tissue). Partial ROIs had either a substantial number of pixels that were white (e.g., if a metallic marker was present) or black (e.g., if the estimated breast border included some non-breast area). From the 108 normal mammograms (two views from each case) there were a total of 20,679 ROIs selected, or approximately 200 per image. From the 10 abnormal cases, 45 ROIs were selected and each ROI contained a portion of the breast cancer that presented as a mass.
4. Each ROI was then reduced in size by averaging 8x8 pixels together. This produced a 15x15 pixel ROI. (One row and one column were lost in MATLAB average subsample algorithm for some unknown reason.) This produced good results visually (see Fig. 2), so no other down sampling methods were tried. Once we are in the training phase, we will try using a median down sampling method to see if we get improved results.

## 5.2. Train classifier to recognize normal mammograms

In our original statement of work, we proposed use a support vector machine (SVM) as our classifier {6}. We have however, decided to use a self-organizing map (SOM) for the following reasons {1}.

1. There are many different appearances of breast lesions (e.g., calcifications, circumscribed masses, spiculated masses, etc.). There are even more different appearances of normal breast

tissues, since the appearance of normal breast tissues depend upon breast thickness, breast density, amount of breast compression, the parenchymal (Wolfe) pattern, position in the breast, etc. Given the wide variety of both normal and abnormal patterns, it would take a very sophisticated (or complex) classifier to class all possible normal and abnormal breast patterns into two classes. SVMs are designed to produce two classes, while SOMs are designed to handle multiple classes.

2. SOM is an unsupervised classifier and SVM is a supervised classifier. The distinction is that for supervised classifiers, one needs to know the classes in the problem. Even if one decided to use multiple classes with a SVM, the classes must be defined *a priori*. However, we do not know *a priori* all the possible different classes. We believe that an unsupervised classifier is ideally suited to this problem, as it will determine the number of classes present in the data.

3. An SVM relies on data that are on the “border” between the two classes. Since most normal patterns are very different from abnormal patterns, any training example that is obviously normal will not be “useful” for training. In this problem, most of the normal training examples will not be useful. An SOM relies on all training samples.

An SOM is useful for reducing multi-dimensional data – 225 (15x15) dimensions in our study – to a two-dimensional surface. An SOM consists of a 2-D array of nodes. Each node represents a category based on a 225-element vector – each element corresponds to one pixel value. This vector is the weights of the SOM. When trained, the SOM adjusts the vector at each node to best match the training data. The first training ROI is compared to each vector at all the nodes. The node that has a vector most similar to the ROI is selected and its vector and those in a neighborhood surrounding the select node are adjust to be more similar to the input ROI. This is repeated for each ROI in the training set, after which one training epoch has been completed. After training, given an input ROI, the SOM will output which node or category that ROI belongs, so the output of the SOM is a number between 1 and the number of nodes. Before training begins, each vector element for each node needs to be initialized. In our study, we used a random number generate to randomly assign values (called weights).

#### 5.2.a. Train Classifier

Since we do not have experience using SOM, we first did some preliminary studies to test the reliability of the SOM for our problem. We tested two different sized SOMs, one was 9x5 (for a total of 45 nodes) and the other was 9x15 (for a total of 135 nodes). The first test was to train the SOM using different number of epochs, to determine the optimum number of training epochs. We tested 100, 200, 400, 800, 1600, and 3200 training epochs. For each training epoch we repeated the training twice. We would expect that after a sufficient number of training epochs, a test ROI would always be placed in the same category and if the SOM is retrained with the same data using different starting weights, the same test ROI should be placed in the same category. Table 1 shows the result for the 9x5 SOM. Clearly, the SOM is not stable. Similar results were obtained for the 9x15 SOM.

We are currently trying to get a stable SOM by changing the learning rate in conjunction with changing the number of epochs and the size of the SOM. We are also trying simpler test problems to further our understanding of SOMs. If these all fail, we will try using less down sampling to preserve more of the pattern of the normal breast. We will also try changing the size of the ROI to larger and smaller sizes. A larger size of ROI will allow the ROI to capture more of the pattern that is present in the image, while a smaller sized ROI will reduce the complexity of the problem for the SOM.



Table 1. Test of the stability of the SOM.

	TEST ROI 1		TEST ROI 2		TEST ROI 3	
Epochs	Run 1	Run 2	Run 1	Run 2	Run 1	Run 2
100	27	14	26	17	13	27
200	23	14	24	14	9	32
400	28	27	23	27	19	9
800	12	29	11	23	27	14
1600	17	27	17	24	27	9
3200	27	14	24	14	16	27
6400	27	27	24	26	13	12

#### 5.2.b. Measure the performance of the technique (using testing dataset)

We have not done this step, since we have not developed a reliable method based on the SOM.

### 5.3 Recommendations in relation to the Statement of Work

We implemented two changes to our original statement of work. First, we preprocessed the images using a wavelet filter instead of two methods proposed: a morphological operator and a linear detection algorithm developed by Zwiggelaar *et al.* This was done because we could not get the latter method to work properly and the wavelet method was faster to implement. Second, we used a self-organizing map (SOM) classifier instead of a support vector machine (SVM). The reason for this change is given in Section 5.2.

## 6. KEY RESEARCH ACCOMPLISHMENTS

- Database of abnormal and normal mammograms has been developed
- Method for reducing image size and preprocessing the images has been developed

## 7. REPORTABLE OUTCOMES

Given the difficulty we have had so far, we do not have any reportable outcomes.

## 8. CONCLUSIONS

We have developed a database and preprocessing method for identifying normal mammograms. The preprocessed database consists of regions-of-interest (ROIs) from normal mammograms and ROIs containing portions of breast cancer from abnormal mammograms. All ROIs have been processed using a wavelet filter to enhance linear structures in the breast.

We are in the process of training a self-organizing map (SOM) to classify the normal and abnormal ROIs. To date, we have not been able to produce a reliable SOM.

Although we have not yet been successful in developing a method to identify normal mammograms using an SOM, the database that we have created can be used to develop other approaches to identifying normal mammograms in the future.

## 9. REFERENCES

1. Kohonen T: *Self-Organizing Maps* (Springer-Verlag, New York, 1997).
2. Nishikawa RM, Giger ML, Schmidt RA, Wolverton DE, and Doi K: Prospective testing of a clinical CAD workstation for the detection of breast lesions on mammograms. In: *Computer Aided Diagnosis in Medical Imaging*. Doi K, MacMahon H, Giger ML, and Hoffmann KR, Eds. (Elsevier, Amsterdam, 1999), pp. 209-214.
3. Bick U, Giger ML, Schmidt RA, Nishikawa RM, Wolverton DE, Lu P, Vyborny CJ, and Doi K: Automated segmentation of digitized mammograms. *Academic Radiology* 2: 1-9, 1995.
4. Strickland RN, and Hahn H: Wavelet transforms for detecting microcalcifications in mammograms. *IEEE Transactions on Medical Imaging* 15: 218-229, 1996.
5. Zwiggelaar R, Parr TC, Schumm JE, Hutt IW, Taylor CJ, Astley SM, and Boggis CR: Model-based detection of spiculated lesions in mammograms. *Med Image Anal* 3: 39-62., 1999.
6. El-Naqa I, Yang Y, Wernick MN, Galatsanos NP, and Nishikawa RM: A support vector machine approach for detection of microcalcifications. *IEEE Transactions on Medical Imaging* 21: 1552-1563, 2002.

## 10. Appendices

None